

## Exhibit A



## Sample GenBank Record

PubMed

Entrez

BLAST

OMIM

Taxonomy

Structure

## GenBank Flat File Format

*Click on any link in this sample record to see a detailed description of that data element or field. All of the descriptions are included on this page, so it can be printed as a single document. You can also return to the Alphabetical Quicklinks Table or Resource Guide*

```

LOCUS      SCU49845      5028 bp      DNA      PLN      21-JUN-1999
DEFINITION Saccharomyces cerevisiae TCP1-beta gene, partial cds, and Axl2p
            (AXL2) and Rev7p (REV7) genes, complete cds.
ACCESSION  U49845
VERSION    U49845.1  GI:1293613
KEYWORDS   .
SOURCE     Saccharomyces cerevisiae (baker's yeast)
  ORGANISM Saccharomyces cerevisiae
            Eukaryota; Fungi; Ascomycota; Saccharomycotina; Saccharomycetes;
            Saccharomycetales; Saccharomycetaceae; Saccharomyces.
REFERENCE  1  (bases 1 to 5028)
  AUTHORS  Torpey,L.E., Gibbs,P.E., Nelson,J. and Lawrence,C.W.
  TITLE    Cloning and sequence of REV7, a gene whose function is required for
            DNA damage-induced mutagenesis in Saccharomyces cerevisiae
  JOURNAL  Yeast 10 (11), 1503-1509 (1994)
  PUBMED   7871890
REFERENCE  2  (bases 1 to 5028)
  AUTHORS  Roemer,T., Madden,K., Chang,J. and Snyder,M.
  TITLE    Selection of axial growth sites in yeast requires Axl2p, a novel
            plasma membrane glycoprotein
  JOURNAL  Genes Dev. 10 (7), 777-793 (1996)
  PUBMED   8846915
REFERENCE  3  (bases 1 to 5028)
  AUTHORS  Roemer,T.
  TITLE    Direct Submission
  JOURNAL  Submitted (22-FEB-1996) Terry Roemer, Biology, Yale University, New
            Haven, CT, USA
FEATURES   Location/Qualifiers
     source          1..5028
                     /organism="Saccharomyces cerevisiae"
                     /db_xref="taxon:4932"
                     /chromosome="IX"
                     /map="9"
     CDS             <1..206
                     /codon_start=3
                     /product="TCP1-beta"
                     /protein_id="AAA98665.1"
                     /db_xref="GI:1293614"
                     /translation="SSIIYNGISTSGLDLNNGTIADMRQLGIVESYKLRVSSASEA
                     AEVLLRVDNIIRARPRTANRQHM"
     gene            687..3158
                     /gene="AXL2"
     CDS             687..3158
                     /gene="AXL2"
                     /note="plasma membrane glycoprotein"
                     /codon_start=1
                     /function="required for axial budding pattern of S.
                     cerevisiae"
                     /product="Axl2p"
                     /protein_id="AAA98666.1"
                     /db_xref="GI:1293615"

```



```


2281 gaagggaaaa tccagacgat gaaaacttac cgcattgctat tagtggaacct gatttgaata
2341 atcctgcaaa taaaccaaat caagaaaacg ctacaccttt gaacaacccc tttgatgatg
2401 atgcttccctc gtacgatgat acttcaatag caagaagatt ggctgctttg aacactttga
2461 aattggataa ccactctgcc actgaatctg atatttccag cgtggatgaa aagagagatt
2521 ctctatcagg tatgaataga tacaatgatc agttccaatc ccaaagtaaa gaagaattat
2581 tagcaaaaacc ccagtagacg cctccagaga gcccgttctt tgaccacacg aataggctct
2641 cttctgtgta tatggatagt gaaccagcag taaataaatc ctggcgatat actggcaacc
2701 tgtcaccagt ctctgatatt gtcagagaca gttacggatc acaaaaaact gttgatacag
2761 aaaaactttt cgatttagaa gcaccagaga aggaaaaacg tacgtcaagg gatgtcacta
2821 tgtcttcact ggacccttgg aacagcaata ttagcccttc tcccgttaaga aaatcagtaa
2881 caccatcacc atataacgta acgaagcatc gtaaccgcca cttacaaaat attcaagact
2941 ctctatccgg taaaaacgga atcactccca caacaatgtc aacttcatct tctgacgatt
3001 ttgttccggt taaagatggt gaaaattttt gctgggtcca tagcatggaa ccagacagaa
3061 gaccaagtaa gaaaaggtta gtagattttt caaataagag taatgtcaat gttggtaag
3121 ttaaggacat tcacggacgc atcccagaaa tgctgtgatt atacgcaacg atattttgct
3181 taattttatt ttctgtttt attttttatt agtggtttac agatacccta tattttattt
3241 agtttttata cttagagaca ttttaatttta attccattct tcaaatttca tttttgact
3301 taaaacaaag atccaaaaat gctctcgccc tcttcatatt gagaatacac tccattcaaa
3361 attttgtcgt caccgtgat taatttttca cttaaactgat gaataatcaa agggccccag
3421 ctcaaacgga ctaaaagaag gagttttatt ttaggaggtt gaaaaccatt attgtctggt
3481 aaattttcat cttcttgaca ttttaaccag tttgaatccc tttcaatttc tgctttttcc
3541 tccaaactat cgaccctcct gtttctgtcc aacttatgtc ctagttccaa ttcgatcgca
3601 ttaataactg cttcaaatgt tattgtgtca tcgttgactt taggtaattt ctccaaatgc
3661 ataatacaac tatttaagga agatcggaat tcgtcgaaca cttcagtttc cgtaatgatc
3721 tgatcgtctt tatccacatg ttgtaattca ctaaaatcta aaacgtattt ttcaatgcat
3781 aaatcgttct ttttattaat aatgcagatg gaaaatctgt aaacgtgcgt taatttagaa
3841 agaactcca gtataagttc ttctatatag tcaattaaag caggatgcct attaatggga
3901 acgaactgcg gcaagttgaa tgactggttaa gtagttagt cgaatgactg aggtgggtat
3961 acattttctat aaaataaaat caaattaatg tagcatttta agtataacct cagccacttc
4021 tctacccatc tattcataaa gctgacgcaa cgattactat ttttttttcc tttctggatc
4081 tcagtcgtcg caaaaacgta taccttcttt ttccgacctt ttttttagct tttctggaaa
4141 gtttatatta gttaaacagg gtctagtctt agtgtgaaag ctagtggttt cgattgactg
4201 atattaagaa agtggaatt aaattagtag tgtagacgta tatgcatatg tatttctcgc
4261 ctgtttatgt ttctacgtac ttttgattta tagcaagggg aaaagaaata catactattt
4321 tttggtaaag gtgaaagcat aatgtaaaag ctagaataaa atggacgaaa taaagagagg
4381 cttagttcat cttttttcca aaaagcacc aatgataata actaaaatga aaaggatttg
4441 ccactctgtc gcaacatcag ttgtgtgagc aataataaaa tcatcacctc cgttgccttt
4501 agcgcgtttg tcgtttgtat cttccgtaat tttagtctta tcaatgggaa tcataaattt
4561 tccaatgaat tagcaatttc gtccaattct ttttgagctt cttcatattt gctttggaat
4621 tcttcgcact tcttttccca ttcactctct tcttcttcca aagcaacgat ccttctacc
4681 atttgctcag agttcaaatc ggctcttttc agtttatcca tgcttctctt cagtttggt
4741 tcaactgtct ctagctgttg ttctagatcc tgggttttct tgggttagtt ctcattatta
4801 gatctcaagt tattggagtc ttcagccaat tgctttgtat cagacaattg actctctaac
4861 ttctccactt cactgtcgag ttgctcgttt ttagcggaca aagatttaac ctctgtttct
4921 ttttcagtg tagattgtc taattctttg agctgttctc tcagctcctc atatttttct
4981 tgccatgact cagattctaa ttttaagcta ttcaatttct ctttgatc

```

//

The corresponding live record for U49845 can be viewed in Entrez.

Examples of other records that show a range of biological features are listed below.

FIELD	COMMENTS
LOCUS	The LOCUS field contains a number of different data elements, including locus name, sequence length, molecule type, GenBank division, and modification date. Each element is described below.
• Locus Name	The locus name in this example is SCU49845. 

The locus name was originally designed to help group entries with similar sequences: the first three characters usually designated the organism; the fourth and fifth characters were used to show other group designations, such as gene product; for segmented entries, the last character was one of a series of sequential integers. (See GenBank release notes section 3.4.4 for more info.)

However, the 10 characters in the locus name are no longer sufficient to represent the amount of information originally intended to be contained in the locus name. The only rule now applied in assigning a locus name is that it must be unique. For example, for GenBank records that have 6-character accessions (e.g., U12345), the locus name is usually the first letter of the genus and species names, followed by the accession number. For 8-character character accessions (e.g., AF123456), the locus name is just the accession number.

The RefSeq database of reference sequences assigns formal locus names to each record, based on gene symbol. RefSeq is separate from the GenBank database, but contains cross-references to corresponding GenBank records.

Entrez Search Field: Accession Number [ACCN]

Search Tip: It is better to search for the actual accession number rather than the locus name, because the accessions are stable and locus names can change.

---

- **Sequence Length**

Number of nucleotide base pairs (or amino acid residues) in the sequence record. In this example, the sequence length is 5028 bp. ↑

There is no maximum limit on the size of a sequence that can be submitted to GenBank. You can submit a whole genome if you have a contiguous piece of sequence from a single molecule type. However, there is a limit of 350 kb on an individual GenBank record (with some exceptions, as noted in section 1.3.2 of the release notes for GenBank 112.0 target="one"). That limit was agreed upon by the international collaborating sequence databases to facilitate handling of sequence data by various software programs. (For more information, see NCBI News articles on Complete Genomes and GenBank Enters Megabase Era.) The minimum length required for submission is 50 bp, although there might be some shorter records from past years.

Entrez Search Field: Sequence Length [SLEN]

Search Tips: (1) To retrieve records within a range of lengths, use the colon as the range operator, e.g., 2500:2600[SLEN]. (2) To retrieve all sequences shorter than a certain number, use 2 as the lower bound, e.g., 2:100[SLEN]. (3) To retrieve all sequences longer than a

certain number, use a series of 9's as the upper bound, e.g., 325000:99999999[SLen].

---

- **Molecule Type**

The type of molecule that was sequenced. In this example, the molecule type is DNA. ↑

Each GenBank record must contain contiguous sequence data from a single molecule type. The various molecule types are described in the Sequin documentation and can include genomic DNA, genomic RNA, precursor RNA, mRNA (cDNA), ribosomal RNA, transfer RNA, small nuclear RNA, and small cytoplasmic RNA.

Entrez Search Field: Properties [PROP]

Search Tip: Search term should be in the format:

**biomol\_genomic**, **biomol\_mRNA**, etc. For more examples, view the Properties field in the Index mode.

---

- **GenBank Division**

The GenBank division to which a record belongs is indicated with a three letter abbreviation. In this example, GenBank division is PLN. ↑

The GenBank database is divided into 18 divisions:

1. PRI - primate sequences
2. ROD - rodent sequences
3. MAM - other mammalian sequences
4. VRT - other vertebrate sequences
5. INV - invertebrate sequences
6. PLN - plant, fungal, and algal sequences
7. BCT - bacterial sequences
8. VRL - viral sequences
9. PHG - bacteriophage sequences
10. SYN - synthetic sequences
11. UNA - unannotated sequences
12. EST - EST sequences (expressed sequence tags)
13. PAT - patent sequences
14. STS - STS sequences (sequence tagged sites)
15. GSS - GSS sequences (genome survey sequences)
16. HTG - HTG sequences (high-throughput genomic sequences)
17. HTC - unfinished high-throughput cDNA sequencing
18. ENV - environmental sampling sequences

Some of the divisions contain sequences from specific groups of organisms, whereas others (EST, GSS, HTG, etc.) contain data generated by specific sequencing technologies from many different organisms. The **organismal divisions are historical and do not reflect the current NCBI Taxonomy**. Instead, they merely serve as a convenient way to divide GenBank into smaller pieces for those who want to FTP the database. Because of this, and because sequences from a particular organism can exist in technology-based divisions such as EST, HTG,

etc., **the NCBI Taxonomy Browser should be used for retrieving all sequences from a particular organism.**

The divisions are also listed in section 3.3 of the GenBank release notes.

The RNA division of GenBank was removed in release 113.0 (August 1999). Sequences that were previously in the RNA division have been moved to the appropriate organismal division. (See section 1.3.2 of the GenBank 113.0 release notes for additional information.)

The HTC division was added to GenBank in release 123.0 (April 2001) and is described in Section 1.3.3 of the GenBank 123.0 release notes.

Another division, called CON, was added in release 115.0 (December 1999) but is not listed above because it records in that division contain no sequence data. Instead, they contain sequence assembly instructions on how to construct contigs from multiple GenBank records. See the Fall 1999 NCBI News and section 1.3.3 of GenBank 115.0 release notes for details.

Entrez Search Field: Properties [PROP]

Search Tip: Search term should be in the format:

**gbdiv\_pri**, **gbdiv\_est**, etc. For more examples, view the Properties field in the Index mode. For example, to eliminate all sequences from a particular division, such as all ESTs, you can use a Boolean query formatted such as:  
human[ORGN] NOT gbdiv\_est[PROP]

For the reasons noted above, **do not use GenBank divisions to retrieve all sequences from a specific organism. Instead, use the NCBI Taxonomy Browser.**

---

- **Modification Date**

The date in the LOCUS field is the **date of last modification**. The sample record shown here was last modified on 21-JUN-1999.



In some cases, the modification date might correspond to the release date, but there is no way to tell just by looking at the record. If you need to know the first date of public availability for a specific sequence record, send a message to [info@ncbi.nlm.nih.gov](mailto:info@ncbi.nlm.nih.gov). We will check the history of the record for you, and let you know the date of first public release. If the sequence was originally submitted to our collaborators at DDBJ or EMBL, rather than to GenBank, we will ask them to send the release date information to you. (See also notes re: date in the Direct Submission reference.)

Entrez Search Field: Modification Date [MDAT]

Search Tips: (1) Enter search term in the format:

yyyy/mm/dd, e.g., 1999/07/25. (2) To retrieve records modified between two dates, use the colon as a range operator, e.g., 1999/07/25:1999/07/31[MDAT]. (3) You can

use the Publication Date [PDAT] field of Entrez to limit search results by the date on which records were added to the Entrez system. Publication date can be in the form of a range, just like the Modification Date.

---

## DEFINITION

Brief description of sequence; includes information such as source organism, gene name/protein name, or some description of the sequence's function (if the sequence is non-coding). If the sequence has a coding region (CDS), description may be followed by a completeness qualifier, such as "complete cds". (See GenBank release notes section 3.4.5 for more info.) ↑

Entrez Search Field: Title Word [TITL]

Search Tip: Although nucleotide definition lines follow a structured format, GenBank does not use a controlled vocabulary, and authors determine the content of their records. Therefore, if a search for a specific term does not retrieve the desired records, try other terms that authors might have used, such as synonyms, full spellings, or abbreviations. The "related records" (or "neighbors") function of Entrez also allows you to broaden your search by retrieving records with similar sequences, regardless of the descriptive terms used by the submitters.

---

## ACCESSION

The unique identifier for a sequence record. An accession number applies to the complete record and is usually a combination of a letter(s) and numbers, such as a single letter followed by five digits (e.g., U12345) or two letters followed by six digits (e.g., AF123456). Some accessions might be longer, depending on the type of sequence record. ↑

Accession numbers do not change, even if information in the record is changed at the author's request. Sometimes, however, an original accession number might become secondary to a newer accession number, if the authors make a new submission that combines previous sequences, or if for some reason a new submission supercedes an earlier record.

Records from the RefSeq database of reference sequences have a different accession number format that begins with two letters followed by an underscore bar and six or more digits, for example:

NT_123456	constructed genomic contigs
NM_123456	mRNAs
NP_123456	proteins
NC_123456	chromosomes

Note: compare accession number with Sequence Identifiers such as Version and GI for nucleotide sequences and protein\_id and GI for amino acid sequences.

Entrez Search Field: Accession [ACCN]

Search Tip: The letters in the accession number can be written in upper- or lowercase. RefSeq accessions must contain an underscore bar between the letters and the numbers, e.g., NM\_002111.

---

## VERSION

A nucleotide sequence identification number that represents a single, specific sequence in the GenBank database. This identification number uses the accession.version format implemented by GenBank/EMBL/DDBJ in February 1999. ↑

If there is any change to the sequence data (even a single base), the version number will be increased, e.g., U12345.1 → U12345.2, but the accession portion will remain stable.

The accession.version system of sequence identifiers runs parallel to the GI number system, i.e., when any change is made to a sequence, it receives a new GI number AND an increase to its version number.

For more information, see section 1.3.2 of the GenBank 111.0 release notes, and section 3.4.7 of the current GenBank release notes.

A Sequence Revision History tool is available to track the various GI numbers, version numbers, and update dates for sequences that appeared in a specific GenBank record (more information and example).

More details about sequence identification numbers and the difference between GI number and version are provided in Sequence Identifiers: A Historical Note.

Entrez Search Field: use the default setting of "All Fields"

---

## • GI

"GenInfo Identifier" sequence identification number, in this case, for the nucleotide sequence. If a sequence changes in any way, a new GI number will be assigned. ↑

A separate GI number is also assigned to each protein translation within a nucleotide sequence record, and a new GI is assigned if the protein translation changes in any way (see below).

GI sequence identifiers run parallel to the new **accession.version** system of sequence identifiers. For more information, see the description of Version, above, and section 3.4.7 of the current GenBank release notes.

A Sequence Revision History tool is available to track the various GI numbers, version numbers, and update dates for sequences that appeared in a specific GenBank record (more information and example).

More details about sequence identification numbers and



the difference between GI number and version are provided in Sequence Identifiers: A Historical Note.

Entrez Search Field: use the default setting of "All Fields"

---

## KEYWORDS

Word or phrase describing the sequence. If no keywords are included in the entry, the field contains only a period. ↑

The Keywords field is present in sequence records primarily for historical reasons, and is not based on a controlled vocabulary. Keywords are generally present in older records. They are **not** included in newer records unless: (1) they are not redundant with any feature, qualifier, or other information present in the record; or (2) the submitter specifically asks for them to be added and #1 is true; or (3) the record contains a special type of sequence such as EST, STS, GSS, HTG, etc.

Entrez Search Field: Keyword [KYWD]

Search Tip: Because keywords are not present in many records, it is best not to search that field. Instead, search All Fields [ALL], the Text Word [WORD] field, or the Title Word [TITL] field, for progressively narrower retrieval.

---

## SOURCE

Free-format information including an abbreviated form of the organism name, sometimes followed by a molecule type. (See section 3.4.10 of the GenBank release notes for more info.) ↑

Entrez Search Field: Organism [ORGN]

Search Tip: For some organisms that have well-established common names, such as baker's yeast, mouse, and human, a search for the common name will yield the same results as a search for the scientific name, e.g., a search for "baker's yeast" in the organism field retrieves the same number of documents as "Saccharomyces cerevisiae". This is true because the Organism field is connected to the NCBI Taxonomy Database, which contains cross-references between common names, scientific names, and synonyms for organisms represented in the Sequence databases.

---

### • Organism

The formal scientific name for the source organism (genus and species, where appropriate) and its lineage, based on the phylogenetic classification scheme used in the NCBI Taxonomy Database. If the complete lineage of an organism is very long, an abbreviated lineage will be shown in the GenBank record and the complete lineage will be available in the Taxonomy Database. (See also the /db\_xref=taxon:nnnn Feature qualifier, below.) ↑

Entrez Search Field: Organism [ORGN]

Search Tip: You can search the Organism field by any

node in the taxonomic hierarchy, e.g., you can search for the term "Saccharomyces cerevisiae", "Saccharomycetales", "Ascomycota", etc. to retrieve all the sequences from organisms in a particular taxon.

---

## REFERENCE

Publications by the authors of the sequence that discuss the data reported in the record. References are automatically sorted within the record based on date of publication, showing the oldest references first.



Some sequences have not been reported in papers and show a status of "unpublished" or "in press". When an accession number and/or sequence data has appeared in print, sequence authors should send the complete citation of the article to [update@ncbi.nlm.nih.gov](mailto:update@ncbi.nlm.nih.gov) and the GenBank staff will revise the record.

Various classes of publication can be present in the References field, including journal article, book chapter, book, thesis/monograph, proceedings chapter, proceedings from a meeting, and patent.

The **last citation** in the REFERENCE field usually contains information about the submitter of the sequence, rather than a literature citation. It is therefore called the "**submitter block**" and shows the words "**Direct Submission**" instead of an article title. Additional information is provided below, under the header **Direct Submission**. Some older records do not contain a submitter block.

Entrez Search Field: The various subfields under References are searchable in the Entrez search fields noted below.

---

### • AUTHORS

List of authors in the order in which they appear in the cited article.



Entrez Search Field: Author [AUTH]

Search Tip: Enter author names in the form: Lastname AB (without periods after the initials). Initials can be omitted. Truncation can also be used to retrieve all names that begin with a character string, e.g., Richards\* or Boguski M\*.

---

### • TITLE

Title of the published work or tentative title of an unpublished work.



Sometimes the words "**Direct Submission**" instead of an article title. This is usually true for the last citation in the REFERENCE field because it tends to contain information about the submitter of the sequence, rather than a literature citation. The last citation is therefore called the "**submitter block**". Additional information is provided

below, under the header **Direct Submission**. Some older records do not contain a submitter block.

Entrez Search Field: Text Word [WORD]

Note: For sequence records, the Title Word [TITL] field of Entrez searches the Definition Line, not the titles of references listed in the record. Therefore, use the Text Word field to search the titles of references (and other text-containing fields).

Search Tip: If a search for a specific term does not retrieve the desired records, try other terms that authors might have used, such synonyms, full spellings, or abbreviations. The 'related records' (or 'neighbors') function of Entrez also allows you to broaden your search by retrieving records with similar sequences, regardless of the descriptive terms used by the submitters.

---

- **JOURNAL**

MEDLINE abbreviation of the journal name. (Full spellings can be obtained from the Entrez Journals Database.)



Entrez Search Field: Journal Name [JOUR]

Search Tip: Journal names can be entered as either the full spelling or the MEDLINE abbreviation. You can search the Journal Name field in the Index mode to see the index for that field, and to select one or more journal names for inclusion in your search.

---

- **PUBMED**

PubMed Identifier (PMID).



References that include PubMed IDs contain links from the sequence record to the corresponding PubMed record. Conversely, PubMed records that contain accession number(s) in the SI (secondary source identifier) field contain links back to the sequence record(s).

Entrez Search Field: It is not possible to search the Nucleotide or Protein sequence databases by PubMed ID. However, you can search the PubMed (literature) database of Entrez for the PubMed ID, and then link to the associated sequence records.

---

- **Direct Submission**

Contact information of the submitter, such as institute/department and postal address. This is always the last citation in the References field. Some older records do not contain the "Direct Submission" reference. However, it is required in all new records.



The Authors subfield contains the submitter name(s), Title contains the words "Direct Submission", and Journal contains the address.

The date in the Journal subfield is the date on which the author prepared the submission. In many cases, it is also the date on which the sequence was received by the

GenBank staff, but it is not the date of first public release. If you need to know the latter, send a message to [info@ncbi.nlm.nih.gov](mailto:info@ncbi.nlm.nih.gov). We will check the history of the record for you.

Entrez Search Field: Use the Author Field [AUTH] if searching for the author name. Use All Fields [ALL] if searching for an element of the author's address (e.g., Yale University). Note, however, that retrieved records might contain the institution name in a field such as Comment, rather than in the Direct Submission reference, so you might get some false hits.

Search Tip: It is sometimes helpful to search for both the full spelling and an abbreviation, e.g., "Washington University" OR "WashU", because the spelling used by authors might vary.

---

## FEATURES

Information about genes and gene products, as well as regions of biological significance reported in the sequence. These can include regions of the sequence that code for proteins and RNA molecules, as well as a number of other features. (See section 3.4.12 of the GenBank release notes for more info.)



A **complete list of features** is available in the following places:

- Appendix III: Feature keys reference of the DDBJ/EMBL/GenBank Feature Table provides definitions, optional qualifiers, and comments for each feature. An alphabetical list is also available.
- Appendix IV: Summary of qualifiers for feature keys provides definitions for the Feature qualifiers.
- Sequin Help documentation (scroll down to 'Features' in the table of contents to see an alphabetical list of features with links to descriptions)
- section 3.4.12.1 of the GenBank release notes

The **location of each feature** is provided as well, and can be a single base, a contiguous span of bases, a joining of sequence spans, and other representations. If a feature is located on the complementary strand, the word "complement" will appear before the base span. If the "<" symbol precedes a base span, the sequence is partial on the 5' end (e.g., CDS <1..206). If the ">" symbol follows a base span, the sequence is partial on the 3' end (e.g., CDS 435..915>).

For more information about feature locations, see the Sequin Help Documentation and section 3.4.12.2 of the GenBank release notes.

The sample record shown here only includes a small number of features (source, CDS, and gene, all of which are described below). The Other Features section, below,

provides links to some GenBank records that show a variety of additional features.

Entrez Search Field: Feature Key [FKEY]

Search Tip: To scroll through the list of available features, view the Feature Key field in Index mode. You can then select one or more features from the index to include in your query. For example, you can limit your search to records that contain both primer\_bind and promoter features.

---

- **source**

Mandatory feature in each record that summarizes the length of the sequence, scientific name of the source organism, and Taxon ID number. Can also include other information such as map location, strain, clone, tissue type, etc., if provided by submitter. ↑

Entrez Search Field: All Fields [ALL] can be used to search for some elements in the source field, such as strain, clone, tissue type.

Use the Sequence Length [SLEN] field to search by length and the Organism [ORGN] field to search by organism name.

Because map location is written as free text and can be represented in a number of ways (e.g., chromosome number, cytogenetic location, marker name, physical map location), it is not directly searchable in the Entrez Nucleotide or Protein databases. However, there are a number of resources that allow you to browse and/or search the maps of various genomes.

---

**Taxon**

A stable unique identification number for the taxon of the source organism. A taxonomy ID number is assigned to each taxon (species, genus, family, etc.) in the NCBI Taxonomy Database. See also the Organism field, above. ↑

Entrez Search Field: The Taxonomy ID number is not seachable in the Organism search field of Entrez but is searchable in the Taxonomy Browser.

Note: The /db\_xref qualifier is one of many that can be applied to various features. A complete list is available in Appendix IV: Summary of qualifiers for feature keys of the DDBJ/EMBL/GenBank Feature Table, and in section 3.4.12.3 of the GenBank release notes. Appendix III: Feature keys reference shows which qualifiers can be used with specific features (see alphabetical list).

---

- **CDS**

Coding sequence; region of nucleotides that corresponds with the sequence of amino acids in a protein (location includes start and stop codons). The CDS feature includes an amino acid translation. Authors can specify the nature of ↑

the CDS by using the qualifier `/evidence=experimental` or `/evidence=not_experimental`.

Submitters are also encouraged to annotate the mRNA feature, which includes the 5' untranslated region (5'UTR), coding sequences (CDS, exon), and 3' untranslated region (3'UTR).

Entrez Search Field: Feature Key [FKEY]

Search Tip: You can use this field to limit your search to records that contain a particular feature, such as CDS. To scroll through the list of available features, view the Feature Key field in Index mode. A complete list of features is also available from the resources noted above.

• **<1..206**

Base span of the biological feature indicated to the left, in this case, a CDS feature. (The CDS feature is described above, and its base span includes the start and stop codons.) **Features can be complete, partial on the 5' end, partial on the 3' end, and/or on the complementary strand.** Examples: ↑

1. **complete** feature is simply written as *n..m*

Example: 687..3158

The feature extends from base 687 through base 3158 in the sequence shown

2. **<** indicates **partial on the 5' end**

Example: <1..206

The feature extends from base 1 through base 206 in the sequence shown, and is partial on the 5' end

3. **>** indicates **partial on the 3' end**

Example: 4821..5028>

The feature extends from base 4821 through base 5028 and is partial on the 3' end

4. **(complement)** indicates that the feature is on the complementary strand

Example: complement(3300..4037)

The feature extends from base 3300 through base 4037 but is actually on the complementary strand. It is therefore read in the opposite direction on the reverse complement sequence. (For an example, see the third CDS feature in the sample record shown on this page. In this case, the amino acid translation is generated by taking the reverse complement of

bases 3300 to 4037 and reading that reverse complement sequence in its 5' to 3' direction.)

---

**protein\_id**

A protein sequence identification number, similar to the Version number of a nucleotide sequence. Protein IDs consist of three letters followed by five digits, a dot, and a version number. If there is any change to the sequence data (even a single amino acid), the version number will be increased, but the accession portion will remain stable (e.g., AAA98665.1 will change to AAA98665.2).

The accession.version format of protein sequence identification numbers was implemented by GenBank/EMBL/DDBJ in February 1999 and runs parallel to the GI number system. More details about sequence identification numbers and the difference between GI number and version are provided in Sequence Identifiers: A Historical Note.

Entrez Search Field: use the default setting of "All Fields"

---

**GI**

"GenInfo Identifier" sequence identification number, in this case, for the protein translation.

The **GI** system of sequence identifiers runs parallel to the **accession.version** system, which was implemented by GenBank, EMBL, and DDBJ in February 1999. Therefore, if the protein sequence changes in any way, it will receive a new GI number, and the suffix of the protein\_id will be incremented by one.

For more information, see the description of protein\_id, above, section 1.3.2 of the GenBank 111.0 release notes, and section 3.4.7 of the current GenBank release notes.

More details about sequence identification numbers and the difference between GI number and version are provided in Sequence Identifiers: A Historical Note.

Entrez Search Field: use the default setting of "All Fields"

---

**translation**

The amino acid translation corresponding to the nucleotide coding sequence (CDS). In many cases, the translations are conceptual. Note that authors can indicate whether the CDS is based on experimental or non-experimental evidence.

Entrez Search Field: It is not possible to search the translation subfield using Entrez. If you want use a string of amino acids as a query to retrieve similar protein sequences, use BLAST instead.

---

- **gene**

A region of biological interest identified as a gene and for which a name has been assigned. The base span for the gene feature is dependent on the furthest 5' and 3' features. Additional examples of records that show the relationship between gene features and other features such as mRNA and CDS are AF165912 and AF090832.

Entrez Search Field: Feature Key [FKEY]

Search Tip: You can use this field to limit your search to records that contain a particular feature, such as a gene. To scroll through the list of available features, view the Feature Key field in Index mode. A complete list of features is also available from the resources noted above.

---

- **complement**

Indicates that the feature is located on the complementary strand.

---

- **Other Features**

Examples of other records that show a variety of biological features; a graphic format is also available for each sequence record and visually represents the annotated features:

- **AF165912** (gene, promoter, TATA signal, mRNA, 5'UTR, CDS, 3'UTR) GenBank flat file
- **AF090832** (protein bind, gene, 5'UTR, mRNA, CDS, 3'UTR) GenBank flat file
- **L00727** (alternatively spliced mRNAs) GenBank flat file

A complete list of features is available from the resources noted above.

---

**ORIGIN**

The ORIGIN may be left blank, may appear as "Unreported," or may give a local pointer to the sequence start, usually involving an experimentally determined restriction cleavage site or the genetic locus (if available). This information is present only in older records.

The sequence data begin on the line immediately below ORIGIN. To view/save the sequence data only, display the record in FASTA format. A description of FASTA format is accessible from the BLAST Web pages.

[Help Desk](#)[NCBI](#)[NLM](#)[NIH](#)[Credits](#)

*Revised October 23, 2006*


*Questions about NCBI resources to* [info@ncbi.nlm.nih.gov](mailto:info@ncbi.nlm.nih.gov)

*Comments about site map to Renata Geer* [renata@ncbi.nlm.nih.gov](mailto:renata@ncbi.nlm.nih.gov)

[Disclaimer](#)   [Privacy statement](#)



## Exhibit B

 Nucleotide banner

My NCBI [\[Sign In\]](#) [\[Register\]](#)

PubMed Nucleotide Protein Genome Structure PMC Taxonomy OMIM Books

Search  for

Limits Preview/Index History Clipboard Details

Display  Show  Send to  Hide: ☐ Sequence ☐ Lesser features

Range: from  to  ☐ Reverse complemented strand Features: ☐ SNP

☐ 1: [AF047436](#). Reports Homo sapiens F1Fo...[gi:3335127]

[Links](#)

[Features](#) [Sequence](#)

LOCUS AF047436 452 bp mRNA linear HTC 21-NOV-2002

DEFINITION Homo sapiens F1Fo-ATPase synthase f subunit mRNA, complete cds.

ACCESSION AF047436

VERSION AF047436.1 GI:3335127

KEYWORDS HTC.

SOURCE Homo sapiens (human)

ORGANISM Homo sapiens  
Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini; Catarrhini; Hominidae; Homo.

REFERENCE 1 (bases 1 to 452)

AUTHORS Mao, M., Fu, G., Wu, J.-S., Zhang, Q.-H., Zhou, J., Kan, L.-X., Huang, Q.-H., He, K.-L., Gu, B.-W., Han, Z.-G., Shen, Y., Gu, J., Yu, Y.-P., Xu, S.-H., Wang, Y., Chen, S.-J. and Chen, Z.

TITLE Identification of genes expressed in human CD34(+) hematopoietic stem/progenitor cells by expressed sequence tags and efficient full-length cDNA cloning

JOURNAL Proc. Natl. Acad. Sci. U.S.A. 95 (14), 8175-8180 (1998)

PUBMED 9653160

REFERENCE 2 (bases 1 to 452)

AUTHORS Zhang, Q.H., Ye, M., Wu, X.Y., Ren, S.X., Zhao, M., Zhao, C.J., Fu, G., Shen, Y., Fan, H.Y., Lu, G., Zhong, M., Xu, X.R., Han, Z.G., Zhang, J.W., Tao, J., Huang, Q.H., Zhou, J., Hu, G.X., Gu, J., Chen, S.J. and Chen, Z.

TITLE Cloning and functional analysis of cDNAs with open reading frames for 300 previously undefined genes expressed in CD34+ hematopoietic stem/progenitor cells

JOURNAL Genome Res. 10 (10), 1546-1560 (2000)

PUBMED 11042152

REFERENCE 3 (bases 1 to 452)

AUTHORS Wu, J.

TITLE Direct Submission

JOURNAL Submitted (10-FEB-1998) Rui-Jin Hospital, Shanghai Second Medical University, Shanghai Institute of Hematology, 197, Rui-Jin Road II, Shanghai 200025, P. R. China

FEATURES

source Location/Qualifiers

1..452

/organism="Homo sapiens"

/mol\_type="mRNA"

/db\_xref="taxon9606"

/cell\_type="CD34+ cell"

CDS

28..312

/codon\_start=1

/product="F1Fo-ATPase synthase f subunit"

/protein\_id="AAC39887.1"

/db\_xref="GI:3335128"

/translation="MASVGECPAPVPVKDKKLLLEVKLGLPSWILMRDFSPSGIFGAFQRGYRYYNKYINVKKSISGITMVLACYVLSYSFSYKHLKHERLRKYH"

ORIGIN

1 gggcacagcg gacaccagga ctccaaaatg gcgtcagttg gtgagtgtcc ggccccagta

61 ccagtgaagg acaagaaact tctggaggtc aaactggggg agctgccaaag ctggatcttg

121 atgcgggact tcagtcctag tggcattttc ggagcgtttc aaagagggtta ctaccggtac

181 tacaacaagt acatcaatgt gaagaagggg agcatctcgg ggattaccat ggtgctggca

241 tgctacgtgc cttttagcta ctctttttcc tacaagcatc tcaagcacga gcggctccgc

301 aaataccact gaagaggaca cactctgcac cccccaccc cagaccttg gcccgagccc

361 ctccgtgagg aacacaatct caatcgttgc tgaatccttt catatcctaa taggaattaa

421 cctccaaata aaacatgact ggtaaaaaaa aa

//

[Disclaimer](#) | [Write to the Help Desk](#)  
[NCBI](#) | [NLM](#) | [NIH](#)

Sep 27 2006 15:22:06



## Sequence Revision History

[PubMed](#)[Nucleotide](#)[Protein](#)[Genome](#)[Structure](#)[PMC](#)[Taxonomy](#)[OMIM](#)[Books](#)Find (Accessions, GI numbers or Fasta style Seqlds) [About Entrez](#) difference between I and II as 

Entrez

Revision history for **AF047436**

Search for Genes  
LocusLink provides curated  
information for human, fruit  
fly, mouse, rat, and  
zebrafish

[Help/FAQ](#)

Batch Entrez: Upload a  
file of GI or accession  
numbers to retrieve  
protein or nucleotide  
sequences

Check sequence  
revision history

How to create WWW  
links to Entrez

[LinkOut](#)[My NCBI \(Cubby\)](#)

Related resources

[BLAST](#)[Reference sequence project](#)[LocusLink](#)[Clusters of orthologous groups](#)[Protein reviews on the web](#)

GI	Version	Update Date	Status	I	II
3335127	1	Nov 21 2002 12:13 PM	Live	<input checked="" type="radio"/>	<input type="radio"/>
3335127	1	May 22 2001 10:38 AM	Dead	<input type="radio"/>	<input checked="" type="radio"/>
3335127	1	Sep 23 1998 4:00 PM	Dead	<input type="radio"/>	<input type="radio"/>
3335127	1	Jul 22 1998 12:15 AM	Dead	<input type="radio"/>	<input type="radio"/>

Accession **AF047436** was first seen at NCBI on Jul 22 1998 12:15 AM

[Disclaimer](#) | [Write to the Help Desk](#)  
[NCBI](#) | [NLM](#) | [NIH](#)

# GenBank

Dennis A. Benson, Ilene Karsch-Mizrachi, David J. Lipman, James Ostell  
and David L. Wheeler\*

National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health,  
Building 38A, 8600 Rockville Pike, Bethesda, MD 20894, USA

Received September 15, 2005; Revised and Accepted October 31, 2005

## ABSTRACT

GenBank (R) is a comprehensive database that contains publicly available DNA sequences for more than 205 000 named organisms, obtained primarily through submissions from individual laboratories and batch submissions from large-scale sequencing projects. Most submissions are made using the Web-based BankIt or standalone Sequin programs and accession numbers are assigned by GenBank staff upon receipt. Daily data exchange with the EMBL Data Library in Europe and the DNA Data Bank of Japan ensures worldwide coverage. GenBank is accessible through NCBI's retrieval system, Entrez, which integrates data from the major DNA and protein sequence databases along with taxonomy, genome, mapping, protein structure and domain information, and the biomedical journal literature via PubMed. BLAST provides sequence similarity searches of GenBank and other sequence databases. Complete bimonthly releases and daily updates of the GenBank database are available by FTP. To access GenBank and its related retrieval and analysis services, go to the NCBI Homepage at [www.ncbi.nlm.nih.gov](http://www.ncbi.nlm.nih.gov).

## INTRODUCTION

GenBank (1) is a comprehensive public database of nucleotide sequences and supporting bibliographic and biological annotation, built and distributed by the National Center for Biotechnology Information (NCBI), a division of the National Library of Medicine (NLM), located on the campus of the US National Institutes of Health (NIH) in Bethesda, MD.

NCBI builds GenBank primarily from the submission of sequence data from authors and from the bulk submission of expressed sequence tag (EST), genome survey sequence

(GSS) and other high-throughput data from sequencing centers. The US Office of Patents and Trademarks also contributes sequences from issued patents. GenBank, the EMBL Data Library (2) in Europe, and the DNA Databank of Japan (DDBJ) (3) comprise the International Nucleotide Sequence Databases and are members of a long-standing collaboration in which data are exchanged daily to ensure a uniform and comprehensive collection of sequence information. NCBI makes the GenBank data available at no cost over the Internet, via FTP and a wide range of Web-based retrieval and analysis services which operate on the GenBank data (4) ([info@ncbi.nlm.nih.gov](mailto:info@ncbi.nlm.nih.gov)).

## ORGANIZATION OF THE DATABASE

From its inception, GenBank has doubled in size about every 18 months. It currently contains over 51 billion nucleotide bases from more than 46 million individual sequences, with 8 million new sequences added in the past year. Contributions from Whole Genome Shotgun (WGS) projects supplement the data in the traditional divisions to bring the total beyond 100 gigabases. Complete genomes ([www.ncbi.nlm.nih.gov/Genomes/index.html](http://www.ncbi.nlm.nih.gov/Genomes/index.html)) represent a growing portion of the database, with over 70 of more than 250 complete microbial genomes in GenBank deposited over the past year. The number of eukaryote genomes for which coverage and assembly are significant continues to increase as well, with over 90 such assemblies now available, including that of the reference human genome.

### Sequence-based taxonomy

Database sequences are classified and can be queried using a comprehensive sequence-based taxonomy ([www.ncbi.nlm.nih.gov/Taxonomy/taxonomyhome.html](http://www.ncbi.nlm.nih.gov/Taxonomy/taxonomyhome.html)) developed by NCBI in collaboration with EMBL and DDBJ and with the valuable assistance of external advisers and curators. Over 205 000 named species are represented in GenBank and new species are being added at the rate of over 3000 per month. About 16% of the sequences in GenBank are of human origin and 11% of

\*To whom correspondence should be addressed. Tel: +1 301 496 2475; Fax: +1 301 480 9241; Email: [wheeler@ncbi.nlm.nih.gov](mailto:wheeler@ncbi.nlm.nih.gov)

© The Author 2006. Published by Oxford University Press. All rights reserved.

The online version of this article has been published under an open access model. Users are entitled to use, reproduce, disseminate, or display the open access version of this article for non-commercial purposes provided that: the original authorship is properly and fully attributed; the Journal and Oxford University Press are attributed as the original place of publication with the correct citation details given; if an article is subsequently reproduced or disseminated not in its entirety but only in part or as a derivative work this must be clearly indicated. For commercial re-use, please contact [journals.permissions@oxfordjournals.org](mailto:journals.permissions@oxfordjournals.org)

all sequences are human ESTs. After *Homo sapiens*, the top species in GenBank in terms of number of bases are *Mus musculus*, *Rattus norvegicus*, *Danio rerio*, *Bos taurus*, *Zea mays*, *Oryza sativa*, *Xenopus tropicalis*, *Canis familiaris* and *Drosophila melanogaster*.

### GenBank records and divisions

Each GenBank entry includes a concise description of the sequence, the scientific name and taxonomy of the source organism, bibliographic references and a table of features ([www.ncbi.nlm.nih.gov/collab/FT/index.html](http://www.ncbi.nlm.nih.gov/collab/FT/index.html)) listing areas of biological significance, such as coding regions and their protein translations, transcription units, repeat regions and sites of mutations or modifications.

The files in the GenBank distribution have traditionally been partitioned into 'divisions' that roughly correspond to taxonomic groups, such as bacteria (BCT), viruses (VRL), primates (PRI) and rodents (ROD). In recent years, divisions have been added to support specific sequencing strategies. These include divisions for EST, GSS, high throughput genomic (HTG), high throughput cDNA (HTC) and environmental sample (ENV) sequences, making a total of 18 divisions. For convenience in file transfer, the larger divisions, such as the EST and PRI, are partitioned into multiple files for the bimonthly GenBank releases on NCBI's FTP site.

### Expressed sequence tags

ESTs continue to be the major source of new sequence records and gene sequences, comprising over 15 billion nucleotide bases in GenBank release 149. Over the past year, the number of ESTs has increased by over 21% to a total of 28.4 million sequences representing more than 740 different organisms. The top five organisms represented in the EST division are *H.sapiens* (6.1 million records), *M.musculus* (4.3 million records), *X.tropicalis* (963 000 records), *Rattus* sp. (701 000 records), *Ciona intestinalis* (684 000 records) and *D.rerio* (635 000 records). As part of its daily processing of GenBank EST data, NCBI identifies through BLAST searches all homologies for new EST sequences and incorporates that information into the companion database, dbEST ([www.ncbi.nlm.nih.gov/dbEST/index.html](http://www.ncbi.nlm.nih.gov/dbEST/index.html)) (5). The data in dbEST are processed further to produce the UniGene database ([www.ncbi.nlm.nih.gov/UniGene/](http://www.ncbi.nlm.nih.gov/UniGene/)) of more than 806 000 gene-oriented sequence clusters representing over 50 organisms, described more fully in (4).

### Sequence-tagged sites (STSs), GSSs and ENV sequences

The STS division of GenBank ([www.ncbi.nlm.nih.gov/dbSTS/index.html](http://www.ncbi.nlm.nih.gov/dbSTS/index.html)) contains over 875 000 sequences, double last year's count, including anonymous STSs based on genomic sequence as well as gene-based STSs derived from the 3' ends of genes and ESTs. These STS records usually include primer sequences, annotations and PCR reaction conditions.

The GSS division of GenBank ([www.ncbi.nlm.nih.gov/dbGSS/index.html](http://www.ncbi.nlm.nih.gov/dbGSS/index.html)) has grown over the past year by 27% to a total of 12.2 million records for over 500 organisms and comprises over 7.6 billion nucleotide bases. GSS records are predominantly single reads from Bacterial Artificial

Chromosomes ('BAC-ends') used in a variety of genome sequencing projects. The most highly represented species in the GSS division are *Zea mays* (1.9 million records), *M.musculus* (1.5 million records), *H.sapiens* (906 000 records) and *C.familiaris* (854 000 records). The human data have been used ([www.ncbi.nlm.nih.gov/genome/clone](http://www.ncbi.nlm.nih.gov/genome/clone)) along with the STS records in tiling the BACs for the Human Genome Project (6).

The ENV division of GenBank, for non-WGS sequences obtained via environmental sampling methods in which the source organism is unknown, debuted with release 147 in April 2005. Records in the ENV division contain 'ENV' in the keyword field and use an 'environmental\_sample' qualifier in the source feature. As of GenBank release 149, the ENV division of GenBank contained over 175 000 sequences, comprising 136 million base pairs, representing more than 3500 studies.

### HTG and HTC sequences

The HTG division of GenBank ([www.ncbi.nlm.nih.gov/HTGS/](http://www.ncbi.nlm.nih.gov/HTGS/)) contains unfinished large-scale genomic records that are in transition to a finished state (7). These records are designated as Phase 0-3 depending on the quality of the data. Upon reaching Phase 3, the finished state, HTG records are moved into the appropriate organism division of GenBank. As of release 149 of GenBank, the HTG division comprised almost 13 billion base pairs of sequence.

The HTC division of GenBank accommodates high-throughput cDNA sequences. HTCs are of draft quality but may contain 5' untranslated regions (5' UTRs) and 3' UTRs, partial coding regions and introns. HTC sequences which are finished and of high quality are moved to the appropriate organism GenBank division. GenBank release 149 contained more than 380 000 HTC sequences totaling over 422 million bases. One project generating HTC data is described in (8).

### Whole genome shotgun sequence

Over 50 million bases of WGS sequence appears in GenBank as sets of WGS contigs, many of them bearing annotations, originating from a single sequencing project. These sequences are issued accession numbers consisting of a 4-letter project ID, followed by a two-digit version number and a 6-digit contig ID. Hence, the WGS accession number 'AAAA 01072744' is assigned to contig number '072744' of the first version of project 'AAAA'. WGS sequencing projects have contributed over 11 million contigs to GenBank, a 3-fold increase over past year's total. These primary sequences have been used to construct some 332 000 large-scale assemblies of scaffolds and chromosomes. WGS project contigs for *H.sapiens*, *C.familiaris*, *Pan troglodytes*, *Drosophila*, *Saccharomyces*, and more than 200 other organisms and environmental samples are available. For a complete list of WGS projects with links to the data, see [www.ncbi.nlm.nih.gov/projects/WGS/WGSprojectlist.cgi](http://www.ncbi.nlm.nih.gov/projects/WGS/WGSprojectlist.cgi).

Submitters of WGS sequences, and genomic sequences in general, are urged to use a new set of evidence tags, described below, in their annotations. In the case of WGS records, these annotations are not tracked from one assembly version to the next and should be considered preliminary.

### New evidence qualifiers

The International Nucleotide Sequence Databases have adopted two new qualifiers to describe the evidence for features annotated in sequence records. The new qualifiers have the form '/experimental=*text*' and '/inference=*TYPE:text*', where '*TYPE*' is one of a number of standard inference types and '*text*' is made up of structured text. These new qualifiers replace 'evidence=experimental' and 'evidence=non-experimental', respectively, which are no longer supported. New versions of the 'tbl2asn' and 'Sequin' sequence submission programs, described below, support the new qualifiers. For details about the new qualifiers and examples of their use, see <http://www.ncbi.nlm.nih.gov/Genbank/evidence.html>.

### BUILDING THE DATABASE

The data in GenBank, and the collaborating databases EMBL and DDBJ, are submitted primarily by individual authors to one of the three databases, or by sequencing centers as batches of EST, STS, GSS, HTC, WGS or HTG sequences. Data are exchanged daily with DDBJ and EMBL so that the daily updates from NCBI servers incorporate the most recently available sequence data from all sources.

### Direct submission

Virtually all records enter GenBank as direct electronic submissions ([www.ncbi.nlm.nih.gov/Genbank/index.html](http://www.ncbi.nlm.nih.gov/Genbank/index.html)), with the majority of authors using the BankIt or Sequin programs. Many journals require authors with sequence data to submit the data to a public database as a condition of publication.

GenBank staff can usually assign an accession number to a sequence submission within two working days of receipt, and do so at a rate of almost 1600 per day. The accession number serves as confirmation that the sequence has been submitted and allows readers of articles in which the sequence is cited to retrieve the data. Direct submissions receive a quality assurance review that includes checks for vector contamination, proper translation of coding regions, correct taxonomy and correct bibliographic citations. A draft of the GenBank record is passed back to the author for review before it enters the database. Authors may ask that their sequences be kept confidential until the time of publication. Since GenBank policy requires that deposited sequence data be made public when the sequence or accession number is published, authors are instructed to inform GenBank staff of the publication date of the article in which the sequence is cited in order to ensure a timely release of the data. Although only the submitting scientist is permitted to modify sequence data or annotations, all users are encouraged to report lags in releasing data or possible errors or omissions to GenBank at [update@ncbi.nlm.nih.gov](mailto:update@ncbi.nlm.nih.gov).

NCBI works closely with sequencing centers to ensure timely incorporation of bulk data into GenBank for public release. GenBank offers special batch procedures for large-scale sequencing groups to facilitate data submission, including the program 'tbl2asn', described at [www.ncbi.nlm.nih.gov/Sequin/table.html](http://www.ncbi.nlm.nih.gov/Sequin/table.html).

### Sequence identifiers and accession numbers

Each GenBank record, consisting of both a sequence and its annotations, is assigned a stable and unique identifier, the accession number, which is shared across the three collaborating databases (GenBank, DDBJ and EMBL) and remains constant over the lifetime of the record even when there is a change to the sequence or annotation. The DNA sequence within a GenBank record is also assigned a unique NCBI identifier, called a 'gi', that appears on the VERSION line of GenBank flatfile records following the accession number. A third identifier of the form 'Accession.version', also displayed on the VERSION line of flatfile records, consolidates the information present in both the gi and accession numbers. An entry appearing in the database for the first time has an 'Accession.version' identifier equivalent to the ACCESSION number of the GenBank record followed by '.1' to indicate the first version of the sequence for the record, e.g. Accession no. AF000001; Version AF000001.1 GI: 987654321.

When a change is made to a sequence given in a GenBank record, a new gi number is issued to the sequence and the version extension of the 'Accession.version' identifier is incremented. The accession number for the record as a whole remains unchanged and the older sequence remains available under the old 'Accession.version' identifier and gi.

A similar system tracks changes in the corresponding protein translations using 'Accession.version' identifiers comprising a protein accession number, e.g. AAA00001, followed by a version number. These identifiers appear as qualifiers for CDS features in the FEATURES portion of a GenBank entry, e.g. /protein\_id='AAA00001.1' Protein sequence translations also receive their own unique gi number, which appears as a second qualifier on the CDS feature, e.g. /db\_xref='GI:1233445'.

### Third Party Annotation (TPA)

TPA records currently support the reporting of published, experimentally confirmed sequence annotation by a scientist other than the original submitter of the primary sequence record in DDBJ/EMBL/GenBank. The scope of the annotations permitted will be expanded in the future to include those derived by inference. TPA sequences may be created by assembling a number of primary sequences. The format of a TPA record (e.g. BK000016) is similar to that of a conventional GenBank record but includes the label 'TPA:' at the beginning of each Definition Line and the keywords 'Third Party Annotation; TPA' in the Keywords field. The Comment field of TPA records lists the primary sequences used to assemble the TPA sequence; the Primary field provides the base ranges of the primary sequences that contribute to the TPA sequence.

Over 4500 TPA records are contained in GenBank release 149, including over 2000 for *D.melanogaster*, 900 for *H.sapiens*, 600 for *O.sativa* and 200 for *M.musculus*. TPA submissions to GenBank may be made using either BankIt, or Sequin but TPA sequences are not released to the public until their accession numbers or sequence data and annotation appear in a peer-reviewed biological journal. For more information on TPA, see [www.ncbi.nlm.nih.gov/Genbank/TPA.html](http://www.ncbi.nlm.nih.gov/Genbank/TPA.html).

### GenBank CON records for assemblies of smaller records

In 1995, the DDBJ/EMBL/GenBank International Nucleotide Sequence Collaboration Databases agreed to a 350 kb limit on the size of most database sequence records in order to conform to the limitations on sequence length of existing molecular biology software. Exceptions were made in the cases of HTG sequence, assemblies of WGS project data and for large eukaryotic genes. The large records that were broken into multiple 350 kb segments to conform to the standard were represented in the GenBank 'CON' division as sets of assembly instructions to allow the transparent display and download of the full record using tools such as NCBI's Entrez. Because of the greater ability of current software to efficiently handle long sequences the 350 kb limit was removed by the Database Collaborators as of June 2004. Although the removal of the limit has immediately allowed many genomes, such as bacterial genomes, to be represented in GenBank as single sequences, it will still be desirable from the standpoints of data transfer and analysis to break some very long sequences, such as portions of eukaryotic genomes, into smaller segments. In these cases, CON division records for the entire sequence will continue to contain assembly instructions to allow the seamless display and download of the sequence.

### BankIt

About one-third of author submissions are received through NCBI's Web-based data submission tool, BankIt ([www.ncbi.nlm.nih.gov/BankIt](http://www.ncbi.nlm.nih.gov/BankIt)). Using BankIt, authors enter sequence information directly into a form and add biological annotation, such as coding regions or mRNA features. Free-form text boxes list boxes, and pull-down menus allow the submitter to further describe the sequence without having to learn formatting rules or restricted vocabularies. BankIt validates submissions, flagging many common errors, and checks for vector contamination using a variant of BLAST called Vecscreen, before creating a draft record in GenBank flat file format for the submitter to review. BankIt is the tool of choice for simple submissions, especially when only one or a small number of records is to be submitted (7). BankIt can also be used by submitters to update their existing GenBank records.

### Sequin and tbl2asn

NCBI also offers a standalone multi-platform submission program called Sequin ([www.ncbi.nlm.nih.gov/Sequin/index.html](http://www.ncbi.nlm.nih.gov/Sequin/index.html)) that can be used interactively with other NCBI sequence retrieval and analysis tools. Sequin handles simple sequences such as a cDNA, as well as segmented entries, phylogenetic studies, population studies, mutation studies, environmental samples, and alignments for which BankIt and other Web-based submission tools are not well-suited. Sequin has convenient editing and complex annotation capabilities and contains a number of built-in validation functions for quality assurance. In addition, Sequin is able to accommodate large sequences, such as that of the 5.6 Mb *Escherichia coli* genome, and read in a full complement of annotations via simple tables. Versions for Macintosh, PC and Unix computers are available via anonymous FTP at <ftp://ftp.ncbi.nih.gov> in the 'sequin' directory. Once a submission is completed, submitters can e-mail the Sequin file to the address [gb-sub@ncbi.nlm.nih.gov](mailto:gb-sub@ncbi.nlm.nih.gov).

Submitters of large, heavily annotated genomes may find it convenient to use 'tbl2asn', referenced above under 'Direct submission', to convert a table of annotations generated via an annotation pipeline, into an ASN.1 record suitable for submission to GenBank.

## RETRIEVING GENBANK DATA

### The ENTREZ system

The sequence records in GenBank are accessible via Entrez ([www.ncbi.nlm.nih.gov/Entrez/](http://www.ncbi.nlm.nih.gov/Entrez/)), a robust and flexible database retrieval system that covers over 30 biological databases containing DNA and protein sequence data, genome mapping data, population sets, phylogenetic sets, environmental sample sets, gene expression data, the NCBI taxonomy, protein domain information, protein structures from the Molecular Modeling Database, MMDB (9), and MEDLINE references via PubMed. The Entrez sequence databases are taken from a variety of sources and therefore include more sequence data than is available within GenBank alone.

### BLAST sequence-similarity searching

Sequence-similarity searches are the most frequent and basic type of analysis performed on the GenBank data. NCBI offers the BLAST ([www.ncbi.nlm.nih.gov/BLAST/](http://www.ncbi.nlm.nih.gov/BLAST/)) family of programs to locate regions of similarity between a query sequence and database sequences (10,11). BLAST searches may be performed on NCBI's Web site, or using a set of stand-alone programs distributed by FTP. BLAST is discussed in a separate article in this issue (4).

### Obtaining GenBank by FTP

NCBI distributes the GenBank releases in the traditional flat-file format as well as in the Abstract Syntax Notation (ASN.1) format used for internal maintenance. The full bimonthly GenBank release and the daily updates, which also incorporate sequence data from EMBL and DDBJ, are available by anonymous FTP from NCBI at (<ftp://ftp.ncbi.nih.gov>) as well as from a mirror site at the University of Indiana (<ftp://bio-mirror.net/biomirror/genbank/>). The full release in flat-file format is available as compressed files in the directory, 'genbank' with a non-cumulative set of updates contained in 'daily-nc'. A script is provided in the 'tools' directory of the GenBank FTP site to convert a set of daily updates into a cumulative update.

### Plans for the future

NCBI has been working with the Consortium for the Barcode of Life (CBOL) ([http://barcoding.si.edu/index\\_detail.htm](http://barcoding.si.edu/index_detail.htm)) to create a new tool for the bulk submission of sequences to GenBank. CBOL is an international initiative devoted to developing DNA barcoding as a tool for characterizing species of organisms using a short DNA sequence derived from a portion of the cytochrome oxidase subunit I gene. Barcode submissions to GenBank include the cytochrome oxidase subunit I sequence combined with a standard set of elements to describe the organism. NCBI offers a new submission tool for Barcode sequences (<http://www.ncbi.nlm.nih.gov/BankIt/barcode/>). Unlike BankIt, which is form based, this web-based submission tool allows users to upload files containing

a batch of sequences with associated source information. It is anticipated that this tool will be used for other types of bulk submissions in the near future.

## CITING GENBANK

If you use the GenBank database in your published research, we ask that this paper be cited.

## ACKNOWLEDGEMENTS

Funding to pay the Open Access publication charges for this article was provided by the National Institutes of Health.

*Conflict of interest statement.* None declared.

## REFERENCES

1. Benson, D.A., Karsch-Mizrachi, I., Lipman, D.J., Ostell, J. and Wheeler, D.L. (2005) GenBank: Update. *Nucleic Acids Res.*, **33**, D34–D38.
2. Kanz, C., Aldebert, P., Althorpe, N., Baker, W., Baldwin, A., Bates, K., Browne, P., van den Broek, A., Castro, M., Cochrane, G. *et al.* (2005) The EMBL Nucleotide Sequence Database. *Nucleic Acids Res.*, **33**, D29–D33.
3. Tateno, Y., Saitou, N., Okubo, K., Sugawara, H. and Gojobori, T. (2005) DDBJ in collaboration with mass-sequencing teams on annotation. *Nucleic Acids Res.*, **33**, D25–D28.
4. Wheeler, D.L., Barrett, T., Benson, D.A., Bryant, S.H., Canese, K., Chetvermin, V., Church, D.M., DiCuccio, M., Edgar, R., Federhen, S. *et al.* (2006) Database resources of the National Center for Biotechnology Information. *Nucleic Acids Res.*, **34**, D173–D180.
5. Boguski, M.S., Lowe, T.M. and Tolstoshev, C.M. (1993) dbEST—database for “expressed sequence tags”. *Nature Genet.*, **4**, 332–333.
6. Smith, M.W., Holmsen, A.L., Wei, Y.H., Peterson, M. and Evans, G.A. (1994) Genomic sequence sampling: a strategy for high resolution sequence-based physical mapping of complex genomes. *Nature Genet.*, **7**, 40–47.
7. Kans, J.A. and Ouellette, B.F.F. (2001) In Baxevanis, A. and Ouellette, B.F.F. (eds), *Bioinformatics: A Practical Guide to the Analysis of Genes and Proteins*. John Wiley and Sons, Inc., New York, NY, pp. 65–81.
8. Hayashizaki, Y., Kawai, J., Shinagawa, A., Shibata, K., Yoshino, M., Itoh, M., Ishii, Y., Arakawa, T., Hara, A., Fukunishi, Y. *et al.* (2001) Functional annotation of a full-length mouse cDNA collection. *Nature*, **409**, 685–690.
9. Marchler-Bauer, A., Anderson, J., Fedorova, N., DeWeese-Scott, C., Geer, L.Y., Hurwitz, D., Jackson, J.J., Jacobs, A., Lanczycki, C., Liebert, C. *et al.* (2005) MMDB: Entrez’s 3D-structure database. *Nucleic Acids Res.*, **33**, D192–D196.
10. Altschul, S.F., Madden, T.L., Schaffer, A.A., Zhang, J., Miller, W. and Lipman, D.J. (1997) Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.*, **25**, 3389–3402.
11. Zhang, Z., Schaffer, A.A., Miller, W., Madden, T.L., Lipman, D.J., Koonin, E.V. and Altschul, S.F. (1998) Protein sequence similarity searches using patterns as seeds. *Nucleic Acids Res.*, **26**, 3986–3991.